

MENU

SEARCH

INDEX

DETAIL

JAPANESE

BACK

4 / 4

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 10-222540

(43)Date of publication of application : 21.08.1998

(51)Int.Cl.

G06F 17/30

(21)Application number : 09-330453

(71)Applicant : N T T DATA TSUSHIN KK

(22)Date of filing : 01.12.1997

(72)Inventor : NAKAJIMA HIROYUKI  
KITANI TSUYOSHI

(30)Priority

Priority number : 08323731

Priority date : 04.12.1996

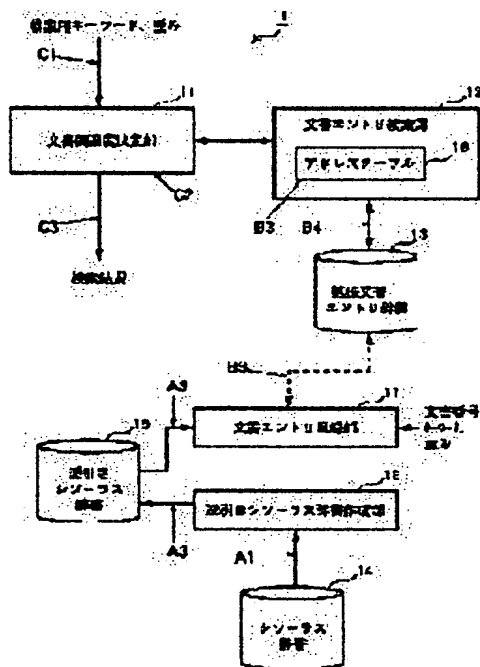
Priority country : JP

## (54) DOCUMENT RETRIEVING METHOD, DEVICE AND RECORDING MEDIUM

(57)Abstract:

**PROBLEM TO BE SOLVED:** To provide a document retrieving device which improves the processing speed of document retrieval.

**SOLUTION:** An extended document entry dictionary 13 stores the document number and weight of a related document in each keyword. The correspondence of a keyword and a document number is preliminarily undergone thesaurus expansion. When a document association degree deciding part 11 inputs a set of a retrieving keyword and weight, a document entry retrieving part 12 reads the document number and weight that correspond to an inputted retrieving keyword from the dictionary 13. The part 11 calculates the degree of association of a retrieving keyword and a document in each read document number and outputs a document number in order of the degree of association.



## LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of rejection]

[Kind of final disposal of application other than the examiner's decision of rejection or application converted registration]

[Date of final disposal for application]



**【特許請求の範囲】**

**【請求項 1】** コンピュータを用いた文書検索方法であって、検索対象となる文書の識別情報、当該文書に含まれるキーワードと検索目的に応じて当該キーワードに付与される重み、及び、前記キーワードが関連する他のキーワードとその関連度を取得して各キーワードに対する文書の重みを算出し、算出された重みを当該文書の識別情報と共にキーワード毎に設定されたメモリ領域へ蓄積する段階と、

検索用キーワードとその重みの入力時に前記検索用キーワードに対応するメモリ領域を参照して当該メモリ領域に存する文書の識別情報とその重みとを抽出し、該抽出結果に基づいて前記検索用キーワードに関連する文書の識別情報を関連度順に出力する段階と、を含むことを特徴とする文書検索方法。

**【請求項 2】** 前記検索対象となる文書に含まれるキーワードと該キーワードが関連する 1 以上の他のキーワードの関連度とを、前記他のキーワードをもとに逆引きできるように予め編集しておき、一のキーワードとその重みが入力されたときに、該入力キーワードが関連するキーワードとその関連度を当該入力キーワードから逆引きして取得し、取得したキーワードとその関連度、及び前記入力キーワードとその重みから、前記各キーワードに対する文書の重みを算出することを特徴とする請求項 1 記載の文書検索方法。

**【請求項 3】** 検索用キーワードと検索目的に応じて当該検索用キーワードに付与された重み、及び、前記検索用キーワードに関連する 1 以上の他のキーワードとその関連度から前記検索用キーワードに対する検索対象となる文書の関連度を決定する文書関連度決定部を有する文書検索装置において、

予め前記文書の識別情報、該文書に含まれるキーワードと検索目的に応じて当該キーワードに付与された重み、及び、前記キーワードが関連する他のキーワードとその関連度に基づいて算出された当該文書の重みを当該文書の識別情報と共にキーワード毎に格納した拡張文書エントリ辞書を備え、

前記文書関連度決定部が、前記文書の関連度を決定する際に、前記拡張文書エントリ辞書から前記検索用キーワードについて格納されている文書の識別情報と当該文書の重みとを抽出することを特徴とする文書検索装置。

**【請求項 4】** 前記文書関連度決定部は、前記検索用キーワードとの関連度の高い順に文書の識別情報を出力することを特徴とする請求項 3 記載の文書検索装置。

**【請求項 5】** 検索対象となる文書に含まれるキーワードと該キーワードが関連する 1 以上の他のキーワードの関連度とを前記他のキーワードをもとに逆引きできるように予め編集されたシソーラス辞書と、一のキーワードと検索目的に応じて前記一のキーワードに付与された重みが入力されたときに、該入力キーワー

ドが関連するキーワードとその関連度を前記シソーラス辞書から逆引きして取得し、取得したキーワードとその関連度、及び前記入力キーワードとその重みから、前記各キーワードに対する文書の重みを算出するとともに、算出された文書の重みと当該文書の識別情報とをキーワード毎に対応付けた拡張文書エントリ辞書を生成する辞書生成手段と、を備え、

検索用キーワードとその重みの入力時に、前記生成された拡張文書エントリ辞書から前記検索用キーワードに対応付けられた文書の識別情報とその重みが抽出されるように構成された文書検索装置。

**【請求項 6】** 前記辞書生成手段は、キーワードとその重み、及び、前記キーワードが関連する他のキーワードとその関連度を取得する度に、各キーワードに対する文書の重みを算出し、算出された重みを当該文書の識別情報と共に前記拡張文書エントリ辞書の同一キーワードの設定領域へ連続的に蓄積することを特徴とする請求項 4 記載の文書検索装置。

**【請求項 7】** 検索対象となる文書の識別情報、当該文書に含まれるキーワードと検索目的に応じて当該キーワードに付与される重み、及び、前記キーワードが関連する他のキーワードとその関連度を取得して各キーワードに対する文書の重みを算出する処理、

算出された重みを当該文書の識別情報と共にキーワード毎に設定された所定のメモリ領域へ蓄積する処理、

検索用キーワードとその重みの入力時に前記検索用キーワードに対応する前記メモリ領域を参照して当該メモリ領域に存する文書の識別情報とその重みとを抽出し、該抽出結果に基づいて前記検索用キーワードに関連する文書の識別情報を関連度順に出力する処理、をコンピュータに実行させるためのプログラムを当該コンピュータが読取可能な形態で記録してなる記録媒体。

**【請求項 8】** 前記プログラムが、さらに、前記検索対象となる文書に含まれるキーワードと該キーワードが関連する 1 以上の他のキーワードの関連度とを、前記他のキーワードをもとに逆引きできるように編集する処理、一のキーワードとその重みが入力されたときに、該入力キーワードが関連するキーワードとその関連度を当該入力キーワードから逆引きして取得し、取得したキーワードとその関連度、及び前記入力キーワードとその重みから、前記各キーワードに対する文書の重みを算出する処理を、コンピュータに実行させることを特徴とする請求項 7 記載の記録媒体。

**【発明の詳細な説明】****【0001】**

**【発明の属する技術分野】** 本発明は、例えば大量の文書や文を蓄積した文書データベース、蓄積された文字等の情報を文書作成や発想展開の支援に用いる各種支援システム等に適用される文書検索技術に係り、特に、入力された検索用キーワードの他に、その検索用キーワードに

関連する他のキーワードをも検索語に加えて該当する文書等を検索する文書検索技術に関する。

#### 【0002】

【従来の技術】 検索対象となる文書や文（以下、文書とする）を格納した文書データベース等から所要の文書を検索するコンピュータ（以下、文書検索装置）において、例えば文書に含まれる検索単語から成るキーワードのみならず、検索目的に応じて付与された重みも入力することで、検索精度を高めることが試みられている。

【0003】 図8は、この種の従来の文書検索装置の構成図である。この文書検索装置8は、個々の文書に含まれるキーワードと当該キーワードに関連する1以上の他のキーワードとをその関連度とともに格納したシソーラス辞書81と、文書識別情報の一例である文書番号とその重みとをキーワード毎に文書エントリとして格納した文書エントリ辞書82とを備えている。

【0004】 また、上記シソーラス辞書81と文書エントリ辞書82とを参照して検索を行うために、シソーラス展開部83、文書関連度決定部84、及び文書エントリ検索部85を備えている。データ量が大きな文書エントリ辞書82については、ハードディスク等の大容量の補助記憶媒体に記録される。

【0005】 シソーラス展開部83は、外部から検索用キーワードとその重み（例えば重要度に応じて付加される数値）の入力を受け付けるとともに、シソーラス辞書81を参照して、検索用キーワードに対応するキーワードを検索し、検索用キーワードと関連キーワードとの関連度に当該検索用キーワードの重みを掛け合わせた関連キーワードの重みを算出する。そして、文書関連度決定部84に検索用キーワードとその重み、及び関連キーワードとその重みを送る。以後、検索用キーワードとそれに対応する関連キーワードとの組をシソーラス展開されたキーワードと称する。

【0006】 文書関連度決定部84は、このシソーラス展開されたキーワードを文書エントリ検索部85に与える。文書エントリ検索部85は、キーワードと文書エントリ辞書82内のメモリ領域のアドレスとの対応関係を記述したアドレステーブル86を保持しており、このアドレステーブル86を参照して、シソーラス展開されたキーワードに対応するアドレスをキーワード毎に取得する。そして、キーワード毎に取得した文書エントリ辞書82のアドレスに存する文書エントリを索出するとともに、それを文書関連度決定部84に返信する。文書関連度決定部84は、文書エントリ検索部85から受け取った文書エントリ（ここでは文書番号）毎に、文書と検索用キーワードとの関連度を算出し、これを検索結果として出力する。

【0007】 ここで、文書と検索用キーワードとの関連度について説明する。前述のように、文書関連度決定部84は、文書エントリ検索部85からシソーラス展開さ

れたキーワード毎に、対応する文書エントリを取得する。これは逆にみれば、文書番号毎に対応する複数のキーワード（シソーラス展開されたキーワード）を取得したことを意味する。そして、文書と検索用キーワードとの関連度は、当該文書番号に対応する一つ以上のキーワードの重みに当該文書番号で表された文書の重みを掛け合わせて得られる数の総和で与えられる。

【0008】 なお、上記関連度の計算、キーワードの重み、文書中のキーワードの重みは、TF/IDF法等のアルゴリズムに基づいて決定される。このTF/IDF法についての詳細については、「Intorduction to Modern Information Retrieval」（Gerald Salton他著、MacGraw-Hill Publishing Company）の記載を参考にすることができる。

#### 【0009】

【発明が解決しようとする課題】 上記従来の文書検索装置8には、シソーラス辞書81に、個々のキーワードについて多くの関連キーワードが格納されている場合、シソーラス展開部83によって検索されるキーワード数が多くなり、これに伴って、文書エントリ検索部85が文書エントリ辞書82において検索するキーワード数が多くなる傾向があった。また、文書エントリ辞書82が大容量の補助記憶媒体に記録されている場合、文書エントリ検索部85が文書エントリ辞書82から文書エントリを検索する際に、適当なシーク時間（アクセスを要求してからアクセス可能になるまでの時間）が必要となる。しかし、一般に大容量の補助記憶媒体のシーク動作は非常に低速であるため、検索速度は、文書エントリ検索部85に渡されるキーワード数の増加に伴って著しく低下する問題があった。

【0010】 そこで、本発明の課題は、キーワードに関連する他のキーワードの数が増大した場合であっても検索速度の低下を防止することができる文書検索方法を提供することにある。本発明の他の課題は、上記文書検索方法の実施に適した文書検索装置、及びこの文書検索方法を汎用のコンピュータを用いて実現するための記録媒体を提供することにある。

#### 【0011】

【課題を解決するための手段】 上記課題を解決する本発明の文書検索方法は、検索対象となる文書の識別情報、当該文書に含まれるキーワードと検索目的に応じて当該キーワードに付与される重み、及び、前記キーワードが関連する他のキーワードとその関連度を取得して各キーワードに対する文書の重みを算出し、算出された重みを当該文書の識別情報と共にキーワード毎に設定されたメモリ領域へ蓄積する段階と、検索用キーワードとその重みの入力時に前記検索用キーワードに対応するメモリ領域を参照して当該メモリ領域に存する文書の識別情報とその重みとを抽出し、該抽出結果に基づいて前記検索用キーワードに関連する文書の識別情報を関連度順に出力

する段階と、を含むことを特徴としている。

【0012】上記方法において、好ましくは、検索対象となる文書に含まれるキーワードと該キーワードが関連する1以上の他のキーワードの関連度とを、前記他のキーワードをもとに逆引きできるように予め編集しておく、一のキーワードとその重みが入力されたときに、該入力キーワードが関連するキーワードとその関連度を当該入力キーワードから逆引きして取得し、取得したキーワードとその関連度、及び前記入力キーワードとその重みから、前記各キーワードに対する文書の重みを算出するようにする。

【0013】上記他の課題を解決するため、本発明は、検索用キーワードと検索目的に応じて当該検索用キーワードに付与された重み、及び、前記検索用キーワードに関連する1以上の他の関連キーワードとその関連度から前記検索用キーワードに対する検索対象となる文書の関連度を決定する文書関連度決定部を有する文書検索装置において、予め前記文書の識別情報、該文書に含まれるキーワードと検索目的に応じて当該キーワードに付与された重み、及び、前記キーワードが関連する他の関連キーワードとその関連度に基づいて算出された当該文書の重みを当該文書の識別情報と共にキーワード毎に格納した拡張文書エントリ辞書を設け、前記文書関連度決定部が、前記文書の関連度を決定する際に、前記拡張文書エントリ辞書から前記検索用キーワードについて格納されている文書の識別情報と当該文書の重みとを抽出し、必要に応じて前記検索用キーワードとの関連度の高い順に文書の識別情報を出力する文書検索装置を提供する。

【0014】本発明は、また、検索対象となる文書に含まれるキーワードと該キーワードが関連する1以上の他のキーワードの関連度とを前記他のキーワードをもとに逆引きできるように予め編集されたシソーラス辞書と、一のキーワードと検索目的に応じて前記一のキーワードに付与された重みが入力されたときに、該入力キーワードが関連するキーワードとその関連度を前記シソーラス辞書から逆引きして取得し、取得したキーワードとその関連度、及び前記入力キーワードとその重みから、前記各キーワードに対する文書の重みを算出するとともに、算出された文書の重みと当該文書の識別情報とをキーワード毎に対応付けた拡張文書エントリ辞書を生成する辞書生成手段と、を備え、検索用キーワードとその重みの入力時に、前記生成された拡張文書エントリ辞書から前記検索用キーワードに対応付けられた文書の識別情報とその重みが抽出されるように構成された文書検索装置をも提供する。

【0015】後者の文書検索装置において、前記辞書生成手段は、例えば、キーワードとその重み、及び、前記キーワードが関連する他のキーワードとその関連度を取得する度に、各キーワードに対する文書の重みを算出し、算出された重みを当該文書の識別情報と共に前記拡張

文書エントリ辞書の同一キーワードの設定領域へ連続的に蓄積するように構成される。

【0016】上記他の課題を解決する本発明の記録媒体は、下記の処理をコンピュータに実行させるためのプログラムがコンピュータ読取可能な形態で記録された記録媒体である。

(1) 検索対象となる文書の識別情報、当該文書に含まれるキーワードと検索目的に応じて当該キーワードに付与される重み、及び、前記キーワードが関連する他のキーワードとその関連度を取得して各キーワードに対する文書の重みを算出する処理、

(2) 算出された重みを当該文書の識別情報と共にキーワード毎に設定された所定のメモリ領域へ蓄積する処理、

(3) 検索用キーワードとその重みの入力時に前記検索用キーワードに対応する前記メモリ領域を参照して当該メモリ領域に存する文書の識別情報とその重みとを抽出し、該抽出結果に基づいて前記検索用キーワードに関連する文書の識別情報を関連度順に出力する処理。

【0017】好ましくは、さらに、下記の処理もコンピュータに実行させるようにする。

(4) 前記検索対象となる文書に含まれるキーワードと該キーワードが関連する1以上の他のキーワードの関連度とを、前記他のキーワードをもとに逆引きできるように編集する処理、

(5) 一のキーワードとその重みが入力されたときに、該入力キーワードが関連するキーワードとその関連度を当該入力キーワードから逆引きして取得し、取得したキーワードとその関連度、及び前記入力キーワードとその重みから、前記各キーワードに対する文書の重みを算出する処理。

#### 【0018】

【発明の実施の形態】以下、図面を参照して、本発明の実施形態を詳細に説明する。図1は、本発明の一実施形態に係る文書検索装置のブロック構成図である。この文書検索装置1は、コンピュータによって実現されるもので、コンピュータの内部あるいは外部記憶装置内に設けられる拡張文書エントリ辞書13、シソーラス辞書14、逆引きシソーラス辞書15、及び、そのコンピュータが所定のプログラムを読み込んで実行することにより形成される、文書関連度決定部11、文書エントリ検索部12、逆引きシソーラス辞書作成部16、文書エントリ登録部17の機能ブロックを備えて構成される。逆引きシソーラス辞書作成部16は、さらに読み出し処理部161、キーワード等登録処理部162の機能要素を含み、文書エントリ登録部17は、検索処理部171、文書エントリ登録処理部172の機能要素を含んでいる。符号173、18はアドレステーブルである。また、図示しないが、装置、検索結果を利用者等に提示するための出力装置をも備えている。

【0019】上記プログラムは、通常、上記内部記憶装置あるいは外部記憶装置に格納され、随時読み取られて実行されるようになっていて、コンピュータとは分離可能な記録媒体、例えばCD-ROMやFD等の可搬性記録媒体、あるいは当該コンピュータ装置と構内ネットワークを通じて接続されたプログラムサーバ等に格納され、使用時に上記内部記憶装置または外部記憶装置にインストールされて随時実行に供されるものであってもよい。

【0020】拡張文書エントリ辞書13には、検索対象となる文書の文書番号と重みとが文書エントリとして格納されており、シソーラス辞書14には、キーワードとそれに関連する1以上の他のキーワード（関連キーワード）がその関連度と共に格納されている。逆引きシソーラス辞書15は、上記シソーラス辞書14の逆引き情報を格納するためのものである。

【0021】文書関連度決定部11は、外部から検索用キーワードとその重みの入力を受け付け、受け付けた検索用キーワードを文書エントリ検索部12に送る。文書エントリ検索部12は、予めキーワードと拡張文書エントリ辞書13の格納領域のアドレスとの対応関係を記述したアドレステーブル18を保持しており、文書関連度決定部11から検索用キーワードを受け取ったとき、このアドレステーブル18を参照して、検索用キーワードに対応する拡張文書エントリ辞書13から該当の文書エントリを読み出し、これを文書関連度決定部11に返信する。文書関連度決定部11は、文書エントリ検索部12から返信された文書エントリ（ここでは文書番号）毎に、文書と検索用キーワードとの関連度を算出し、これを検索結果として出力する。

【0022】ところで、文書検索装置1において上記文書検索を行う場合は、拡張文書エントリ辞書13への文書エントリの登録が完了していることが必要となる。そのため、文書検索装置1では、シソーラス辞書14を参照して逆引きシソーラス辞書作成部16で逆引きシソーラス辞書15を作成しておく。そして、文書エントリ登録部17で、入力されたキーワード及びその重みを上記逆引きシソーラス辞書15を用いて拡張文書エントリ辞書13へ登録する。

【0023】まず、逆引きシソーラス辞書作成部16における処理を図2及び図3を参照して具体的に説明する。逆引きシソーラス辞書作成部16は、図2に示すように、読み出し処理部161においてシソーラス辞書14の内容を読み出す。図3（a）は、ここで読み出されるシソーラス辞書14の内容例を示す図表A1であり、キーワード「kwd1」に関連する関連キーワードが「kwd4」と「kwd7」で、「kwd1」との関連度がそれぞれ“0.8”と“0.4”であることを示している。なお、図3（a）において、キーワード項目はハッシュキーであり、関連キーワード項目はハッシュ結

果を示すものである。他の関連キーワードについても同様な対応関係が与えられる。

【0024】読み出し処理部161は、シソーラス辞書14から読み出した上記内容に基づいて、関連キーワードからキーワードへの逆の対応関係（逆引きの関係）を与える。図3（b）は、このような対応関係の一例を示す図表A2である。つまり、図3（a）によれば、シソーラス辞書14において、キーワード「kwd1」には関連度が“0.8”の関連キーワード「kwd4」と、関連度が“0.4”の関連キーワード「kwd7」とが対応していた。これを逆にみると、関連キーワード「kwd4」は関連度“0.8”のキーワード「kwd1」が対応し、関連キーワード「kwd7」は関連度“0.4”のキーワード「kwd1」に対応したものとなっている。図3（b）は、こうした逆の対応関係を示している。

【0025】なお、図3（b）では、一つの関連キーワードに一つのキーワードを対応させている。しかし、一つの関連キーワード「kwd4」に上記キーワード「kwd1」のほかに他のキーワード「kwd10」も対応する場合がある。そこで、キーワード登録処理部162は、このような対応関係を整理して、一の関連キーワードに複数のキーワードが対応する場合は、キーワード項目に複数のキーワードとその関連度とを格納する。この図表A3から、関連キーワード「kwd4」は、“関連度”0.8”でキーワード「kwd1」と対応し、関連度“0.2”でキーワード「kwd10」に対応していることがわかる。

【0026】次に、逆引きシソーラス辞書15が作成された後の文書エントリ登録部17による登録処理を図4及び図5を参照して説明する。文書エントリ登録部17は、図4に示すように、登録対象となる文書の文書番号と、当該文書についてTF/IDF法等により、あるいは他の手法によって求めたキーワード及び重みとの組を検索処理部171に入力する。図5（a）は、入力されたキーワード及びその重みの例を示す図表B11、文書番号の例を示す図表B12である。図示の例では、文書番号として“120”、この文書番号“120”の文書に含まれるキーワード「kwd4」、「kwd5」の重みが、それぞれ“3”、“2”となっている。

【0027】検索処理部171は、この入力キーワード「kwd4」、「kwd5」に対応するキーワードを逆引きシソーラス辞書15から索出し、各キーワード毎の重みを算出する。このとき、各キーワードの重みは、入力キーワードについてはその重み、索出されたキーワードについてはキーワード間の関連度に入力キーワードの重みを掛け合わせたものである。

【0028】例えば、逆引きシソーラス辞書15の内容が図3（c）で示されたものであるときの検索処理部171の検索処理の結果は、図5（b）の図表B2のよう

になる。つまり、図3(c)によれば、入力キーワード「kwd4」に対応するキーワードは「kwd1」, 「kwd10」である。そして、各キーワード「kwd1」, 「kwd10」の関連度がそれぞれ“0.8”, “0.2”であるので、これらに入力キーワード「kwd4」の重み“3”を掛け合わせて、“2.4”, “0.6”となる。同様にして、他の入力キーワード「kwd5」についても、キーワード「kwd5」の重み“2”と、これに関連するキーワード「kwd2」, 「kwd1」の重み“0.8”, “0.2”を算出する。

【0029】文書エントリ登録部172は、キーワードと拡張文書エントリ辞書13の格納領域のアドレスの対応を示したアドレステーブル173を参照して、キーワード毎に、文書番号と上記重みを登録する。図5(c)はアドレステーブル173の内容例を示す図表B3である。

【0030】このテーブル173では、キーワード「kwd1」, 「kwd2」, …に、それぞれ拡張文書エントリ辞書13の先頭アドレス“1”, “4000”, …が対応付けられている。つまり、アドレス“1”~“3999”がキーワード「kwd1」の領域、“4000”~“”がキーワード「kwd2」の領域、…である。文書エントリ登録部172は、図5(d)の図表B4に示されるように、キーワード「kwd1」の領域に文書番号“120”と重み“2.4”を文書エントリとして登録し、キーワード「kwd2」の領域に文書番号“120”と重み“0.8”を文書エントリとして登録する。他のキーワード「kwd3」, …についても同様にして登録を済ませる。

【0031】なお、ここでは関連キーワードからキーワードの検索が可能な逆引きシソーラス辞書15を使用した。既存のシソーラス辞書14を、さらに関連キーワードからキーワード検索が可能になるように作成すれば、逆引きシソーラス辞書15の代わりに両方向からの検索が可能に作成されたシソーラス辞書14を用いることができる。

【0032】次に、図6及び図7を参照して、上記文書検索装置1を用いた文書検索方法について説明する。文書検索装置1は、文書関連度決定部11において検索用キーワードとその重みの入力を受け付け(ステップS101)、検索用キーワードについては、それを文書エントリ検索部12に送る。図7(a)は入力された検索用キーワードと重みの一例を示す図表C1である。ここには、3つの検索用キーワード「kwd1」, 「kwd2」, 「kwd5」と、各検索用キーワードの重み“1”, “2”, “5”が与えられている。この場合、検索されるのは、これらの3つの検索用キーワードとの関連度が総合的に高い文書(その識別情報)である。

【0033】文書エントリ検索部12は、アドレステー

ブル18を参照して、検索用キーワードに対応する拡張文書エントリ辞書13のアドレスからすでに登録された文書エントリを読み出す(ステップS102)。そして、読み出した文書エントリを文書関連度決定部11に返信する。ここで、アドレステーブル18の内容は、すでに図5(c)で紹介したアドレステーブル173の内容を示す図表B3と同一とし、拡張文書エントリ辞書13の内容は、図5(d)に示した図表B4と同一とする。このとき、図7(a)に示された検索用キーワードが文書エントリ検索部12に送られると、文書エントリ検索部12は、アドレステーブル18を参照して、拡張文書エントリ辞書15におけるキーワード「kwd1」の領域から、そこに登録された文書エントリ、すなわち、文書番号“24”, “120”, “12”, …とそれぞれ文書に与えられた重み“3.4”, “2.4”, “1.2”を読み出す。

【0034】キーワード「kwd2」についても同様に、拡張文書エントリ辞書15におけるキーワード「kwd2」の領域から、そこに登録された文書番号“1002”, “64”, “120”, …とそれぞれ文書に与えられた重み“2.4”, “1.2”, “0.8”を読み出す。キーワード「kwd3」についても同様の処理を行う。図7(b)は、このようにして文書関連度決定部11が文書エントリ検索部12を介して取得したキーワードとその重み、そして、検索された文書の文書番号、文書の重みとの対応関係を示した図表C2である。

【0035】文書関連度決定部11は、文書エントリ検索部12からキーワード毎の文書エントリを受け取った後、その文書エントリで表された文書番号毎に、文書と検索用キーワードとの関連度を算出し(ステップS103)、関連度の大きなものから順に文書番号を関連度とともに検索結果として出力する(ステップS104)。ステップS104での処理を図7(b)の内容を例に挙げて説明する。

【0036】いま、文書番号“24”の文書に注目したとき、重みが“1”のキーワード「kwd1」における文書番号“24”の文書の重みは“3.4”、重みが“2”のキーワード「kwd2」における文書番号“24”の文書の重みは“1.3”、そして重みが“5”のキーワード「kwd3」における文書番号“24”の文書の重みは“0.8”である。こうしてキーワード「kwd1」からの寄与“3.4”、キーワード「kwd2」からの寄与“2.6”、キーワード「kwd1」からの寄与“4.0”を合算すると全体で“10.0”となる。これが文書番号“24”の文書の検索用キーワードに対する関連度である。他の文書についても同様にして、文書と検索用キーワードとの関連度が算出される。

【0037】図7(c)は、文書関連度決定部11の検索結果例を示す図表C3である。ここでは、文書番号“24”の文書が最も大きな関連度“10.0”をも

ち、文書番号“12”の文書が次に大きな関連度“8.9”をもつことを示している。以上が文書検索装置1の検索処理についての説明である。

【0038】拡張文書エントリ辞書13がハードディスク等の補助記憶媒体に格納された場合、文書エントリ検索部12がこの補助記憶媒体にアクセスする際に、適当なシーク時間を必要とする。そして、検索所要時間は、おおよそこのシーク時間の合計で決まる。この事実注目し、本実施形態による文書検索装置1の方が従来技術による文書検索装置8より検索速度において向上していることを説明する。

【0039】なお、説明を簡単にするため、シソーラス辞書81、14には、それぞれキーワードとして唯一のキーワード「kwd」と、それに関連する他のキーワードとしてn個のキーワード「kwd1」～「kwdn」があるものとする。また、文書エントリ辞書81にはキーワード「kwd」、「kwd1」～「kwdn」のそれぞれについて文書エントリがm個登録され、拡張文書エントリ辞書13にはキーワード「kwd」に対して文書エントリが $(n+1) \cdot m$ 個登録されているものとする。さらに文書エントリ辞書81と拡張文書エントリ辞書13に対する文書エントリ検索部85、12のシーク時間をいずれも $t_1$ とし、文書エントリ辞書81と拡張文書エントリ辞書13から一つの文書エントリを読み出すのに必要な時間をいずれも $t_2$ とする。

【0040】このとき、キーワード「kwd」を検索用キーワードとしたとき、従来の文書検索装置8において検索に要する時間を算出する。この場合、シソーラス展開部83から文書関連度決定部84に渡されるキーワード数(=キーワード「kwd」についてシソーラス展開したときのキーワードの数)は、 $n+1$ 個である。この $n+1$ 個のそれぞれについて、文書エントリ辞書81からm個の文書番号を読み出すことになるが、一つのキーワードについて必要なシーク時間は $t_1 + m t_2$ なので、シーク時間の合計は、 $(n+1) \cdot (t_1 + m t_2)$ (=T1)となる。これが文書検索装置7の検索所要時間である。

【0041】一方、本実施形態の文書検索装置1における検索所要時間を算出する。この場合、文書関連度決定部11から文書エントリ検索部12に渡されるキーワード「kwd」は一個である。従って、このキーワード「kwd」について、シーク時間の合計は $t_1 + (n+1) \cdot m t_2$ (=T2)となる。こうして $T_1 - T_2 = n t_1 > 0$ が導かれる。

【0042】なお、本実施形態では、シソーラス辞書81とシソーラス辞書14にキーワードとして唯一のキーワード「kwd」、それに関連する関連キーワードとして「kwd1」～「kwdn」があるような簡単な場合を説明したが、より複雑な場合にも同様な計算を行うことで、本実施形態による文書検索装置1が文書検索装置

8より検索速度において向上していることを確認することができる。

【0043】また、拡張文書エントリ辞書13は、文書エントリ辞書81と比べて登録された文書エントリが増大するが、各文書エントリをキーワード毎の領域に連続して蓄積(登録)しているので、文書エントリの取得時に補助記憶媒体から連続的にこれを読み出すことができ、実質的な速度の低下は防止されている。また、拡張文書エントリ辞書13から一つの文書エントリを取得する際のシーク動作は、文書エントリ辞書81に対するシーク時間と等しいから、シーク動作が少ない分だけ、処理速度を短縮させることができる。

【0044】また、本実施形態では、文書エントリ検索部12及び文書エントリ登録処理部172に、それぞれ図5(c)の図表B3に示した内容のアドレステーブル18、173を備えた例を示したが、このアドレステーブル18、173は、文書エントリ検索部12や文書エントリ処理部172ではなく、拡張文書エントリ辞書13に共通に備えておくこともできる。あるいはアドレステーブル18、173自体を省略することもできる。

【0045】さらに、本発明は、文の検索にも応用することができる。このときには、拡張文書エントリ辞書15に格納する際に、文書番号の代わりに文を識別するための番号を使用すればよい。

【0046】

【発明の効果】以上の説明から明らかなように、本発明によれば、シソーラス展開したキーワード数が増大した場合にも、従来のように改めてシソーラス展開する必要がないので、文書の検索所要速度を向上させることができる効果がある。

【図面の簡単な説明】

【図1】本発明の一実施形態による文書検索装置のブロック構成図。

【図2】逆引きシソーラス辞書作成部の詳細構成図。

【図3】(a)～(c)は逆引きシソーラス辞書の作成過程での具体例を示す図表。

【図4】文書エントリ登録部の詳細構成図。

【図5】(a)～(d)は文書エントリ登録処理過程での具体例を示す図表。

【図6】文書検索装置を用いた文書検索方法の手順説明図。

【図7】(a)～(c)は本実施形態による文書検索の際の具体例を示す図表。

【図8】従来の文書検索装置のブロック構成図。

【符号の説明】

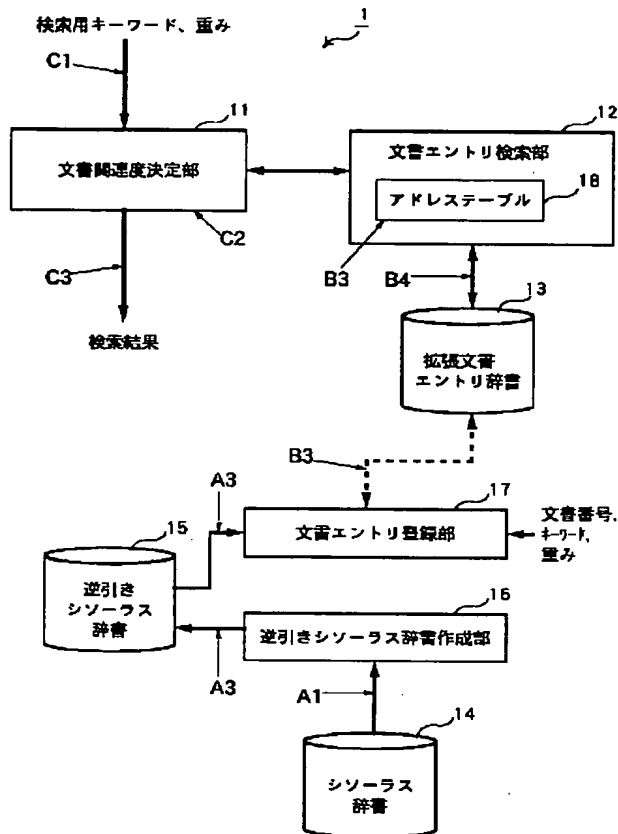
- 1、8 文書検索装置
- 11、84 文書関連度決定部
- 12 文書エントリ検索部
- 13 拡張文書エントリ辞書
- 14、81 シソーラス辞書



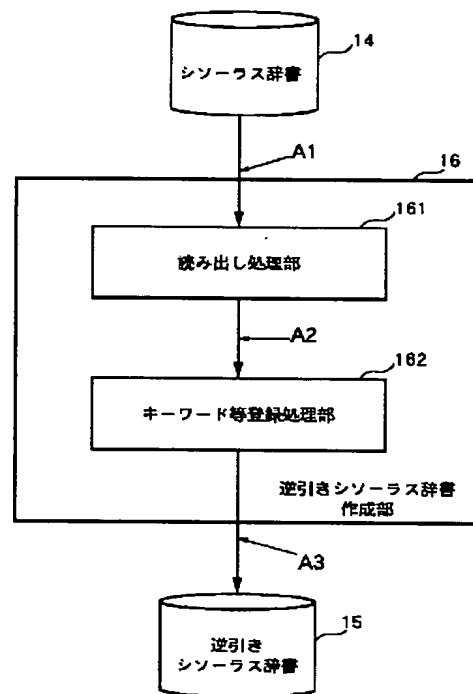
- 15 逆引きシソーラス辞書
- 16 逆引きシソーラス辞書作成部
- 17 文書エントリ登録部
- 82 文書エントリ辞書
- 83 シソーラス展開部
- 85 文書エントリ検索部

- 161 読み出し処理部
- 162 キーワード等登録処理部
- 171 検索処理部
- 172 文書エントリ登録処理部
- 18、173 アドレステーブル

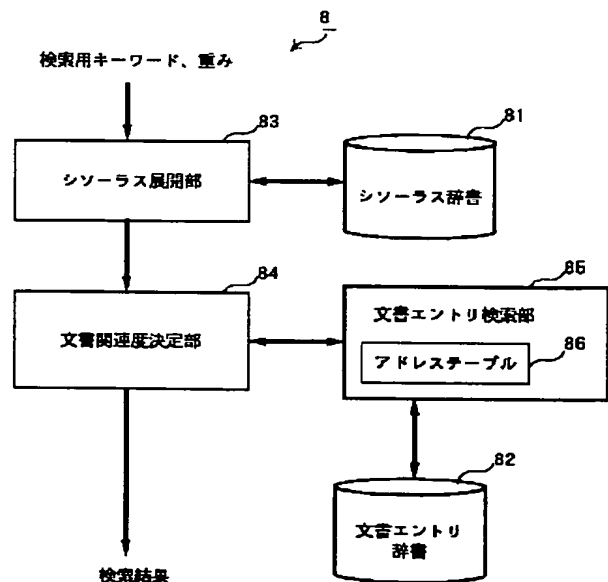
【図 1】



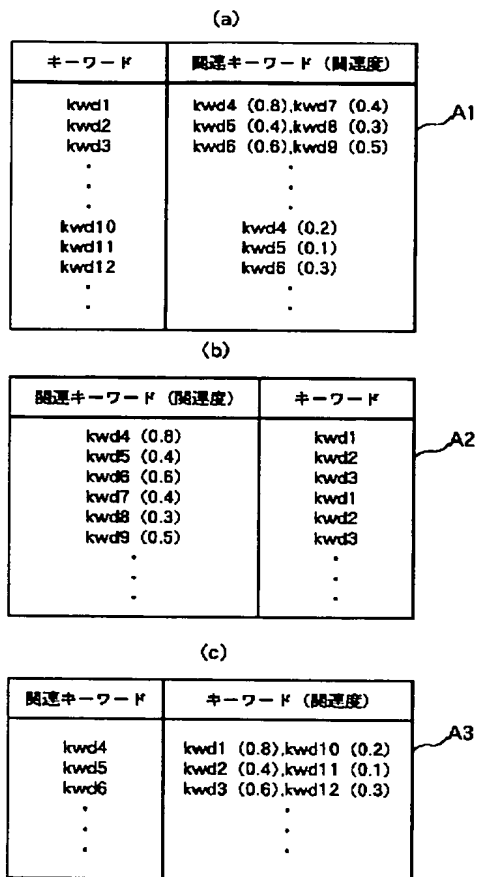
【図 2】



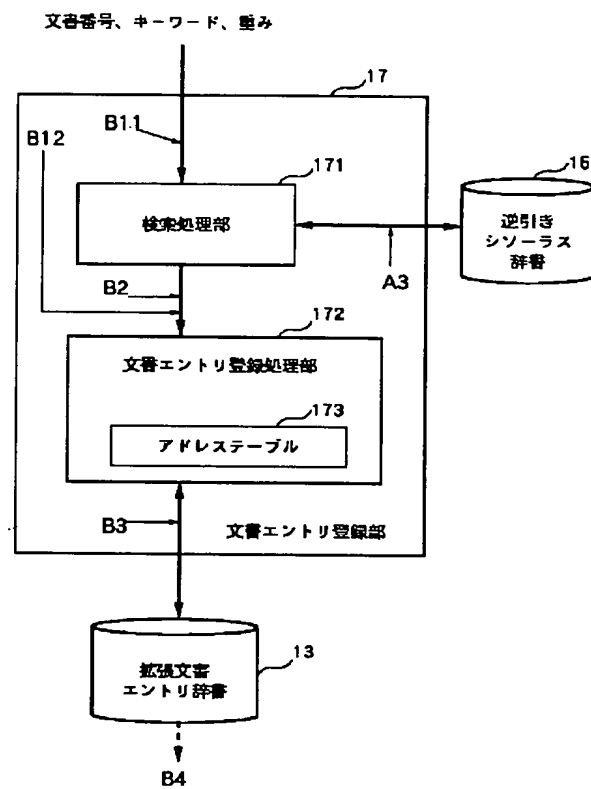
【図 8】



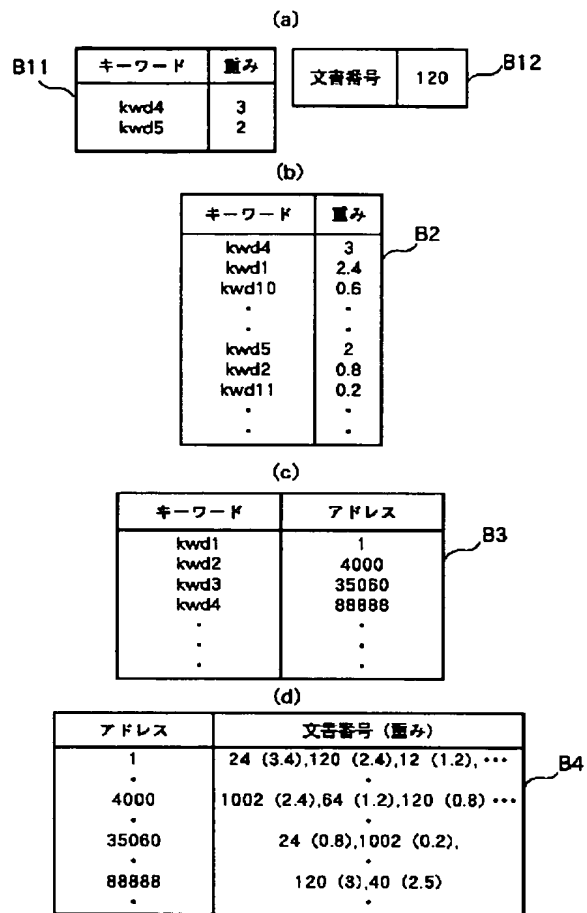
【図 3】



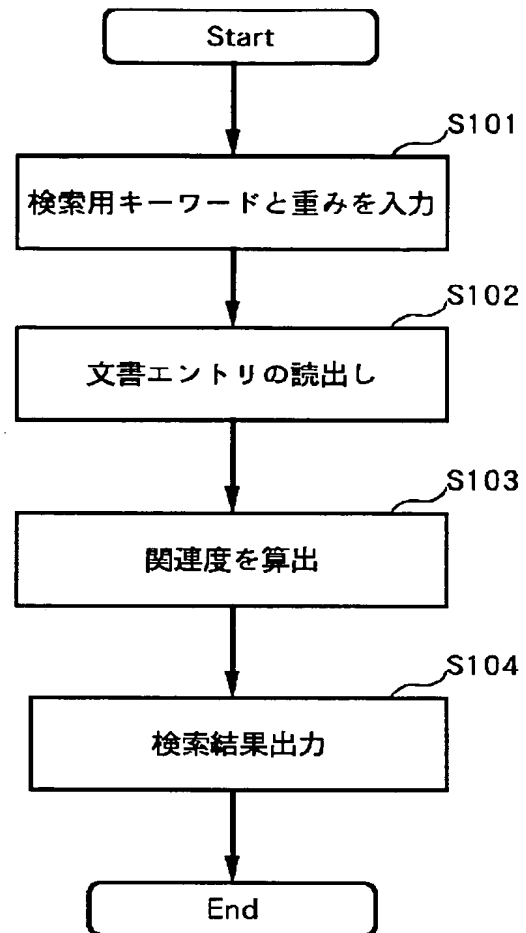
【図 4】



【図 5】



【図 6】



【図 7】

